

# Attentional Push: Augmenting Salience with Shared Attention Modeling

Siavash Gorji     James J. Clark

Centre for Intelligent Machines, Department of Electrical and Computer Engineering,  
McGill University  
Montreal, Quebec, Canada  
siagorji@cim.mcgill.ca clark@cim.mcgill.ca

**Abstract.** We present a novel visual attention tracking technique based on Shared Attention modeling. Our proposed method models the viewer as a participant in the activity occurring in the scene. We go beyond image saliency and instead of only computing the power of an image region to pull attention to it, we also consider the strength with which other regions of the image push attention to the region in question. We use the term *Attentional Push* to refer to the power of image regions to direct and manipulate the attention allocation of the viewer. An attention model is presented that incorporates the Attentional Push cues with standard image saliency-based attention modeling algorithms to improve the ability to predict where viewers will fixate. Experimental evaluation validates significant improvements in predicting viewers' fixations using the proposed methodology in both static and dynamic imagery.

**Keywords:** Visual Attention, Shared Attention, Image Saliency

## 1 Introduction

Attention is a temporal selection mechanism in which a subset of available sensory information is chosen for further processing. Since the visual system cannot perform all visual functions at all locations in the visual field at the same time in parallel [1], attention implements a serialized mechanism that acts as an information-processing bottleneck to allow near real-time performance. Given the wider arrangement of receptors and the larger receptive fields of ganglion cells in the periphery, attention supports analysis of a scene by successively directing the high-resolution fovea to salient regions of the visual field. While visual attention guides the so called *focus of attention* (FOA) to important parts of the scene, a key question is on the computational mechanisms underlying this guidance. Aside from being an interesting scientific challenge, attention tracking-determining where, and to what, people are paying attention while viewing static photographs or while watching videos and cinematic movies- has many applications in: object object detection and recognition [2], visual surveillance [3], human-robot interaction [4], and advertising [5].

Modeling visual attention has attracted much interest recently and there are several frameworks and computational approaches available. The current

state-of-the-art of attention prediction techniques are based on computing image saliency maps, which provide, for each pixel, its probability to attract viewers' attention. Almost all attention models are directly or indirectly inspired by cognitive findings. The basis of many attention models dates back to Treisman and Gelade's feature integration theory [6] which showed that during visual perception, visual features, e.g. color, size, orientation, direction of movement, brightness and spatial frequency, are registered early, automatically, and in parallel across the whole visual field. Koch and Ullman [7] proposed a feed-forward neural model to combine these early visual features into a central representation, i.e. the saliency map. Clark and Ferrier [8] developed a robotic vision system that used the Koch and Ullman saliency model to control the motion of a binocular pair of cameras. This work was the first to demonstrate computationally the link between image saliency and eye movements. Subsequently, models of saliency have often been characterized by how well they predict eye movements.

Perhaps the first complete implementation of the Koch and Ullman model was proposed by the pioneering work of Itti et al. [9] which inspired many later models and has been the standard benchmark for comparison. This model generates feature maps across different scales for three early visual features and then linearly combines them to yield the saliency map. Similarly, GBVS [10] extracts intensity, color, and orientation feature maps at multiple scales and builds a fully connected graph over all locations of each feature map, with weights between two nodes set proportional to the similarity of feature values and their spatial distance. The saliency map is formed by a normalized combination of the equilibrium distribution of the graphs. Goferman et al. [11] proposed a context-aware saliency detection model. The method is based on four principles of human attention: local low-level features such as color and contrast, global considerations to maintain features that deviate from the norm, visual organization rules, and high-level factors such as human faces. In RARE [12], the saliency map is formed by fusing rarity maps, which are computed using cross-scale occurrence probability of each pixel. In AWS [13], the local variability in energy is used as an estimation of saliency. The method decomposes the a and b color channels into multiple scales, while decomposing the luminance channel using Gabor filter banks. The saliency map is computed as the local average of the decomposed channels. In BMS [14], an image is characterized by a set of binary images, generated by randomly thresholding the image's color channels. Based on a Gestalt principle of figure ground segregation, the method computes the saliency map using the topological structure of Boolean maps.

The above models only rely on bottom-up influences. While having reasonable performance, bottom-up models are mostly feed-forward, do not need training and are in general easy to apply. While many attention models fall into this category, they cannot fully explain the eye movements, since the fixations are also modulated by the visual tasks. In contrast to bottom-up attention, top-down attention is slow, task-driven, voluntary, uses feedback and requires learning mechanisms to be trained for a specific visual task and are therefore, more complex to deploy. Top-down attention takes higher-level cognitive cues such as

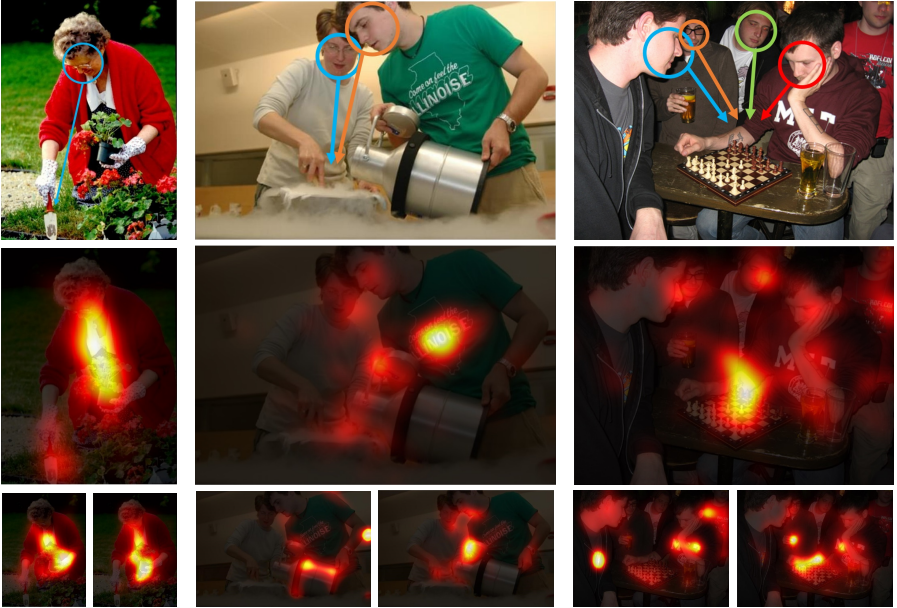


Fig. 1: Static gaze direction as an Attentional Push cue, directing viewers’ attention in social scenes. Each image has a shared locus of attention of the scene actors that has low saliency, in spite of having viewers’ attention allocated to them. (Top row) Original images with annotated head pose. (Middle row) Overlaid fixation maps. (Bottom row) Overlaid saliency maps: (left) BMS [22], (right) eDN [23]. The saliency maps cannot fully predict viewers’ fixations. Original images and eye fixation data are from the action and the social categories of CAT2000 dataset [24]). Saliency maps were histogram-matched to the fixation maps for visualization.

task demands into account. This is probably why regardless of the important role of top-down factors in directing visual attention, the majority of existing attention models focus on bottom-up cues (see the recent extensive survey of attention modeling by Borji and Itti [15]). Haji-Abolhassani and Clark [16] developed an inverse Yarbus process in which the attention tracking system is able to infer the viewer’s visual task, given the eye movement trajectories. Similar methods were proposed by Borji and Itti [17] using a Boosted Classifier and by Kanan et al. [18] using a Fisher Kernel Learning method. Aside from the visual task demands, scene gist [19], tendency of observers to look near the center of displays (also known as image center-bias [20]), and expertise with similar scenes [21], also affect attention in a top-down manner.

All of the aforementioned methods are based on saliency maps, and only differ in their choice of features to be used in forming the maps, and in the way top-down guidance modulates the saliency. In a recent comparative study, Borji et al. [25] compared 35 state-of-the-art of saliency models over synthetic and natural stimuli. They showed that these methods are far from completely predicting viewers’ attentional behavior. A possible reason for this mediocre performance is that image saliency is not the only factor driving attention allocation.

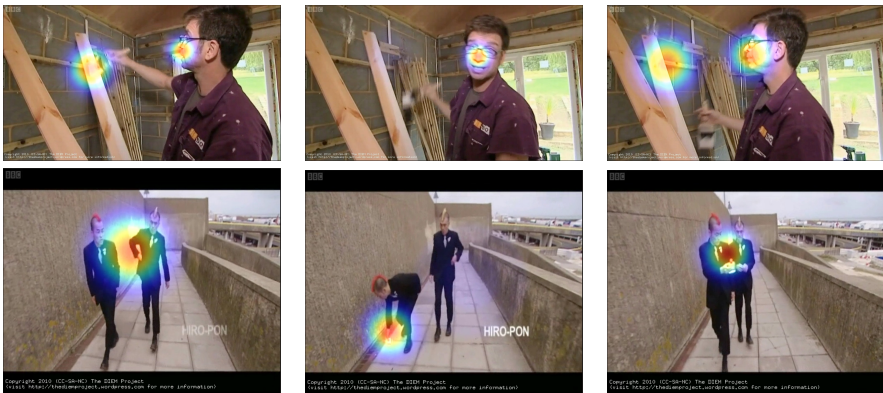


Fig. 2: Dynamic changes in gaze and body pose direction as Attentional Push cues. The images show overlaid fixation maps for three video frames, before, while and after a dynamic gaze/body change. In all cases, viewers’ attention is highly influenced by the Attentional Push cues. Original video and eye fixation data are from the DIEM dataset [28]).

Birmingham et al. [26] assessed the ability of the Itti et al. [9] saliency map in predicting eye fixations in social scenes and showed that its performance is near chance levels. They concluded that the viewer’s eye movements are affected by their interest to social information of the scenes. In a recent study, Borji et al. [27] investigated the effect of gaze direction on the bottom-up saliency. They conducted a controlled experiment in which an actor is asked to look at two different objects in turn, resulting in two images that differed only by the actor’s gaze direction. The experiments show that the median of the fraction of all saccades that start from the head and end inside the gazed-at object to that of the ignored object is more than 3. This clearly shows that low-level saliency cannot account for the influence of gaze direction on fixations. The study also highlights that the median of the saccade directions in the actor’s gaze direction is about 9.5 times higher than the chance level, which indicates that viewers tend to look more in the direction of actor’s gaze than in other directions [27].

One of the shortcomings of the current approaches is that, for the most part, they concentrate on analyzing regions of the image for their power to attract attention. However, as noted above, in many instances, a region of the image may have low saliency, but nonetheless still have attention allocated to it. Clearly, in such cases there are no salient features that attract attention to these regions. Instead, we propose that a viewer has their attention pushed to these regions by some high level process. This suggests that in building an attention model we should go beyond image saliency and instead of only computing the power of an image region to *pull* attention to it, we should also consider the strength with which other regions of the image *push* attention to the region in question.

Our proposed method models the viewer as a passive participant in the activity occurring in the scene. While the viewers cannot affect what is going on in the scene, their attentional state can nonetheless be influenced by the actors

in the scene. We will treat every image viewing situation as one of *Shared Attention*, which is the process by which multiple agents mutually estimate, direct and follow each others attentional state [29]. As one of the building blocks for social communication, shared attention is a bilateral process by which an agent attends to an object that another agent attends to. Here, an agent may refer to both a scene actor and the viewer. To achieve shared attention, agents must observe, coordinate and influence their behaviors in order to engage in a collaborative intentional action [29].

We use the term *Attentional Push* [30] to refer to the power of image regions to direct and manipulate the attention allocation of the viewer. Attentional Push can arise from many sources, which are mostly abstract high-level features, such as faces and body pose. For example in Fig. 1, the head pose and the body pose of the scene actors manipulate the attention of the viewer. Such Attentional Push cues direct the viewers’ attention to the shared locus of attention of the scene actors. Fig. 1 shows that although the shared loci of attention might have low saliency, they have viewers’ attention allocated to them nonetheless. It is also clear that two of the best-performing saliency methods (according to the MIT saliency benchmark [31]), BMS [22] and eDN [23], perform poorly in predicting the fixation maps for such images with social clues. In addition, the strength by which an Attentional Push cue directs the viewers’ attention could intensify as more actors focus their attention to the same shared locus of attention.

We propose that the effect of Attentional Push in directing viewers’ attention intensifies in more immersive scenarios, such as dynamic videos, 3-D movies and ultimately, while using virtual reality setups. Therefore, comparing to standard image saliency-based methods, the prediction performance of an Attentional Push-based method would become more noticeable, as viewers feel more immersed in the ongoing event in the scene. Fig. 2 illustrates the effect of dynamic changes in gaze and body pose direction, as Attentional Push cues, on viewers’ attention, while watching a dynamic movie. It suggests that as the level of immersion increased, viewers’ attention is more influenced by Attentional Push cues.

This paper presents an attention tracking method that combines Attentional Push cues with standard image saliency-based algorithms to improve the ability to predict where viewers’ fixations in social scenes. Our approach to Shared Attention is to first identify the actors in the image, which can then be analyzed for their Attentional Push, potentially directing and manipulating the attention allocation of the viewer. The introduction of attention tracking and prediction techniques based on treating the viewer as a participant in a shared attention situation will open new avenues for research in the attention field.

In a recent study, Parks et al. [32] proposed the DWOC model, an attention model which combines bottom-up saliency with the head pose of the scene actors. The method is based on a two-state Markov chain describing the transition probabilities between head region and non-head region states, which is used to predict whether the next fixation is gaze related or being saliency driven. Our proposed method differs from Parks et al. [32] in the following aspects: (i)

their method only considers the effect of actors’ head pose in manipulating the viewer’s attention, whereas our Shared Attention-based method generalizes to all such Attentional Push cues; (ii) their method is only applicable to static scenes, whereas our method explicitly benefits from dynamic Attentional Push cues in directing viewers’ attention while watching dynamic imagery; (iii) their method requires the viewers’ eye movements to predict the next fixation point, whereas our method is based the image information only; and (iv) their method assumes the viewers have to fixate upon the head regions, in order for their next fixations to be influenced by the actors’ gaze direction. However, this might not be the case and in our model the viewers’ attention might be affected when the viewer tries to understand the gist of the scene.

The rest of this paper is organized as follows. Section 2 elaborates using Attentional Push in attention tracking. Section 3 presents our attention tracking model which augments standard saliency maps with Attentional Push cues. Section 4 illustrates experimental evaluation of the proposed method. Section 5 concludes the paper.

## 2 Attentional Push

To benefit from the Attentional Push cues in predicting viewers’ attention, we propose to consider the viewer of the imagery as a partner in a shared attention situation, where the other partner(s) are the actors in the imagery. The goal of an agent in a shared attention setting is to coordinate its attention with other agents. To achieve this, the agent may try to interpret the intentions of another agent by watching its movements and its attentional behavior. While Kaplan and Hafner [29] require the both agents to be able to detect, manipulate, coordinate and understand the attentional state and the behavior of the other agent in order to reach shared attention, our particular situation is a restricted asymmetric form of shared attention, in that the viewer has no control over the attentional state of the actors in the imagery. However, the actors in the image are assumed to have some control over the attentional state of the other actors in the image, as well as that of the viewer. Our working assumption will be that if two or more actors in a scene have a shared attentional locus, then the viewer will also be compelled to direct his or her attention to that locus. Thus, not only are we tracking the attention of the viewer, we are also tracking the attention of the actors in the scene, and doing so in a cooperative manner.

Many Attentional Push cues have been reported in the literature of attention tracking. Perhaps the most prominent of these are gaze cues. Development of gaze following capabilities for robots via different learning mechanisms has been in the spotlight of research into socially interactive robots human-robot interaction (see the recent survey by Ferreira and Dias [4] and the references therein). Castelhana et al. [33] showed that while the actor’s face is highly likely to be fixated, the viewer’s next saccade is more likely to be toward the object that is fixated by the actor, compared to any other direction. Ricciardelli et al. [34] showed that perceived gaze enhances attention if it is in agreement with the

task direction, and inhibits it otherwise. They showed that in spite of top-down knowledge of its lack of usefulness, the perceived gaze automatically acts as an attentional cue and directs the viewer’s attention. Similarly, as illustrated in Fig. 2, the body pose of the scene actors could also push the viewers’ attention. Although the attentional manipulation strength of the gaze direction dominates the body pose direction in most cases, it could be still intensified if the body pose direction is in agreement with the gaze direction.

Apart from gaze and body pose cues, one of the most frequently cited Attentional Push cues in the literature is the center bias. Borji et al. [25] showed that a simple 2D Gaussian shape drawn at the center of the image predicts the viewers’ fixations well. We can treat the center-bias effect in the shared attention setting by considering the photographer as an actor in the shared attention setting, which tries to put the semantically interesting and therefore, salient elements in the center of the frame. In [35], Tseng et al. showed that center bias is strongly correlated with photographer bias, rather than the viewing strategy and motor bias. There are some attention tracking models (e.g. Judd et al. [36]) that have explicitly used the center-bias as a location prior to achieve better performance in predicting the eye movements.

Aside from the static Attentional Push cues mentioned above, Attentional Push cues can also arise from dynamic events. For example, Smith [37] showed that sudden movements of the heads of actors are a very strong cue for attention, where the viewer’s FOA is not the head itself, but where it is pointing to (see Fig. 2). Smith [37] also notes the ”bounce” in the attention of a movie viewer back to the center of the movie screen when tracking an object which moves off the screen to one side. Similarly, in [35], abrupt scene changes are used to assess the contribution of the center bias in predicting viewer’s attention while watching dynamic stimuli. We believe that employing such Attentional Push cues, either in static or in dynamic scenes, along with bottom-up image saliency would be necessary to predict viewer’s eye movements.

### 3 Augmented Saliency

In this section, we present our attention tracking method which fuses the Attentional Push and the standard image saliency techniques into a single attention tracking scheme. The proposed approach provides a framework for predicting viewer’s FOA while watching static or dynamic imagery. For the sake of readability, the model focuses upon the interaction between one actor and the viewer, although this can be readily adapted in the case of multiple actors by providing unique identifiers for each actor. Our model distinguishes between two sets of attentional cues: Attentional Push-based and saliency-based, and provides a selection mechanism between them. While the saliency-based cues represent properties of the scene objects, the Attentional Push cues are based on the scene actor(s), such as head pose, body pose and dynamic changes in any of them as well as rapid scene changes. The need for a deterministic selection mechanism stems from the fact that in certain circumstances, an Attentional Push cue might

pull the viewer’s attention. An example of such situation is when a scene actor has frontal head directions. This traditional signal of Attentional Pull strictly pulls the viewer’s attention to the actor rather than pushing it elsewhere. This has been exploited in many researches on gaze imitation and Shared Attention (e.g. see [38] and [39]). In the top row of Fig. 2, it could be seen that while the actor’s head pose pushes the viewers’ attention when the actor is looking sideways, it pulls the viewers’ attention when the head pose is frontal. Therefore, it is vital to have a selection mechanism between pulling and pushing viewers’ attention.

Assuming that the scene is observable via an image  $I$ , we can model the actor’s attentional focus  $A$  as conditionally dependent on the bottom-up factors such as location and appearance properties of the scene objects  $\mathbf{O} = \{O_1, \dots, O_k\}$ , as well as the top-down factors of the ongoing task of the scene, parameterized by  $\mathbf{T}$ . We can then describe the attentional manipulation of the scene actors and the scene objects over the viewers’ attention  $V$  by employing a set of latent attentional cues  $\{a_i\}$ . In this Shared Attention setting, the attentional focus of the scene actors and the viewers are given by  $P(A|\mathbf{O}, \mathbf{T})$  and  $P(V|\{a_i\})$ , respectively. Learning and inferring the viewers’ attention using the above dependencies requires the attentional foci of the scene actors. However, in most cases, the eye movements of the scene actors are not available. We hypothesize that this is not actually needed and we can directly employ some overt attentional measures of the actors, such as head gaze direction, body pose direction and hand gesture direction, to infer the viewer’s attention.

As shown in Figure 3, we model the dependency between the attentional focus of the scene actors and the viewers by a set of  $n$  observable Attentional Push cues  $\mathbf{s} = \{s_i^b, s_i^g\}$  and similarly, we use a set of Attentional Pull cues  $\{l_i\}_{i=1}^m$ , arising from image saliency. The graphical model is used as a convenient method to describe the conditional dependencies of Attentional Push-based and saliency-based cues. We employ normalized saliency maps  $S(I)$  to estimate the joint distribution over the set of Attentional Pull cues  $P(l_1, \dots, l_m, l|I)$ . We represent each Attentional Push cue using two distinct quantities: 1) a geometrical structure  $\mathbf{g} : \{x, y, \boldsymbol{\theta}, r, \sigma\}$ , describing the  $(x, y)$  location, 3-D rotation angles ( $\boldsymbol{\theta} = \{\text{roll}, \text{pitch}, \text{yaw}\}$ ) (for symmetrical Attentional Push cues,  $\boldsymbol{\theta}$  is set to the frontal direction), scale ( $\sigma$ ) and confidence factor ( $r$ ); and 2) a variable  $b$  representing the presence or absence of the cue. For static Attentional Push cues,  $b \in [0, 1]$ , while for dynamic Attentional Push cues, we encode the habituation factor [40], i.e. the strength or probability of the viewers’ motor response to a certain stimulus, by  $b(t) := b(0)e^{-\beta(t-t_0)}$ , where  $\beta$  denotes the decay rate,  $t_0$  is the moment of occurrence in which  $b(0)$  is set and  $t$  denotes discretized frame time.

We encode the deterministic constraints of the attentional guidance in the push-pull control node  $C$  in Figure 3. This node’s value is deterministically assigned by its parents, using a predefined set of rules. For each Attentional Push cue  $s_i$ , we construct a 2-D Attentional Push map  $M(s_i)$ , having the advantage of being directly comparable with saliency maps. For directional Attentional Push



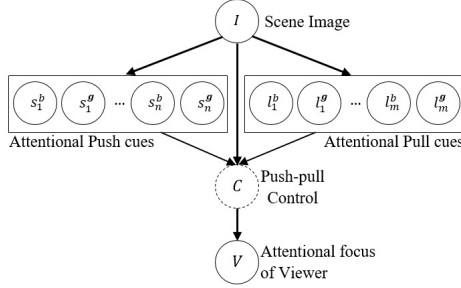


Fig. 3: Shared Attention-based augmented saliency. The viewer’s FOA  $V$  is affected by the set of  $n$  Attentional Push cues, represented by  $\mathbf{s} = \{s_i^b, s_i^g\}$ , and the set of  $m$  Attentional Pull cues, denoted by  $\mathbf{l} = \{l_i\}$ . The model assumes the attentional cues to be directly observable via the scene image  $I$ . The deterministic node  $C$  controls the transitions between the Attentional Push and Attentional Pull cues, based on their current observed values.

cues, i.e. head/body pose and dynamic changes in head/body pose, we represent a 2-D map, having 1s along the direction  $\phi_i$  ( $\phi_i$  denotes the projection of  $\theta_i$  on the image plane), modulated by a 1-D Gaussian function centered at each point with a standard deviation proportional to  $\sigma_i$  in the direction perpendicular to  $\phi_i$  by  $N(s_i)$ . For symmetrical Attentional Push cues, i.e. frontal head pose, center bias, attentional bounce and rapid scene changes, we denote by  $G(s_i)$  a 2-D map, containing a symmetric 2-D Gaussian, centered at the center of the 2-D map, with unit variance. The control node computes the Attentional Push maps by combining the directional and the symmetrical maps as follows:

$$M(s_i) = b(t)[\alpha\sigma_i G(s_i) + (1 - \alpha)N(s_i)]. \quad (1)$$

where  $\alpha$  is 1, if  $\theta_i$  is near frontal and 0 otherwise.

We employ the fusion mechanism in [41] to combine the Attentional Push and Attentional Pull cues by assigning deterministic weights to each of them using their relevant statistics. For Attentional Pull cues  $\mathbf{l} = \{l_i\}$ , we use the mean absolute skewness  $\gamma$ , i.e. the average of the absolute value of the third moments, of the normalized saliency map and for each Attentional Push map  $\{M(s_i)\}$ , we use its confidence measure  $r_i$  in computing the weights. The control node output is determined by

$$C(\mathbf{s}, \mathbf{l}, I) = \gamma S(I) + \sum_{i=1}^n r_i M(s_i) + \gamma S(I) \sum_{i=1}^n r_i M(s_i). \quad (2)$$

Note that the third term in 2, the element-wise multiplication of the saliency map and each Attentional Push map, acknowledges the fact that the directional Attentional Push maps are not able to discern between any image regions in the pose direction. The element-wise multiplication enables the directional Attentional Push-based cues and the saliency-based cues to interact in a way that if both of them have large values on a region, that region would have high saliency in the augmented saliency map.

## 4 Evaluation and Comparison

### 4.1 Estimating Attentional cues

To evaluate the performance of the Attentional Push-based method in predicting viewers' fixations, we employ the following Attentional Push cues: actors' body and head pose, the central bias, changes in actors' head and body pose, the bounce of attention and rapid scene changes. To identify the scene actors, we proceed by detecting humans and faces in the scene. To detect humans, we employ the HoG-based detector of Dalal and Triggs [42]. To detect faces, we use the face detection system of Viola and Jones [43] and deformable mixture of parts-based method of Zhu and Ramanan [44]. Our experiments showed that the combination of the above methods results in a better detection rate, while increasing the false positive rate. For dynamic scenes, the scene actors might have non-frontal head poses which causes most face detection algorithms to fail. Therefore, we employ the state-of-the-art tracker TLD [45], comprising of a median flow-based tracker, a detector, to localize the appearance of the faces, and a learning component which estimates the detector's error and updates it. The method returns a bounding box, computed from the merged results of the tracker and the detector. If neither the tracker nor the detector return a bounding box, the face is declared as non-visible which triggers a bounce of attention cue. To estimate the head pose of the scene actors and their dynamic changes, we employ facial landmarks detection algorithms to accurately estimate the roll, pitch and yaw angles of the actor's head. Here, we use the iterative approach of [46] which initializes the landmarks locations using the face bounding box and uses an incremental cascaded linear regression to update the landmarks locations. To estimate the body pose direction, we use the poselet-based method of Maji, Bourdev and Malik [47]. To detect rapid scene changes, we adopt the method in [48] which is based on comparing the edge strength and orientation of consecutive video frames.

### 4.2 Evaluation protocol

Attention models have commonly been validated against eye movements of human observers. To evaluate the proposed method, we employed three popular image and video datasets: 1) The CAT2000 dataset [24], 2) the NUSEF dataset [49], and 3) The DIEM dataset [28], containing eye movement data from 250 subjects watching 85 different dynamic scenes such as movie trailers, sport events and advertisements. Since the proposed Attentional Push-based method requires actors' in the scene, for the static stimuli, we used all the available images from the Action and the Social categories of the CAT2000 dataset (200 images in total). We also use 150 images from the NUSEF dataset. The employed images (350 images in total) contain humans and faces with resolution high enough for successful detection and accurate pose estimation. Note that if we run the proposed method for images with no actors, the results would be the same as the employed saliency method. For the dynamic stimuli, we use 13 videos from the

DIEM dataset that contain people interacting with each other, each containing more than 1000 video frames (14109 video frames in total). We compare our Attentional Push-based augmented saliency method with the ten best-performing state-of-the-art saliency models, according to the MIT saliency benchmark [31] (see Table 1). For each saliency method, we create an augmented saliency using the proposed methodology. To evaluate attention models, many evaluation metrics have been proposed in the literature (e.g. [31,25]). However, the performance of a model may change remarkably when different metrics are used. To ensure that the main qualitative conclusions are independent of the choice of metric, we analyze the performance of the proposed model using three popular evaluation metrics: the Area Under the ROC Curve (*AUC*), the Normalized Scan-path Saliency (*NSS*), and the Correlation Coefficient (*CC*). To compute *AUC*, fixated points are considered as the positive set while other locations are randomly sampled to form a negative set. By applying multiple thresholds, the saliency map is used as a binary classifier and its ROC curve is plotted as the true positive rate against the false positive rate. Perfect prediction leads to an *AUC* value of 1.0, while random prediction has an *AUC* of 0.5. The *NSS* metric uses the average value of the saliency map, normalized to zero mean and unit variance, at fixation locations. When  $NSS \geq 1$ , the saliency map exhibits significantly higher saliency values at human fixated locations compared to other locations. The *CC* metric measures the strength of a linear relationship between the saliency map and the fixation map. Value of  $abs(CC)$  close to 1 show a perfect linear relationship.

### 4.3 Results and Discussion

Table 1 compares the prediction performance of the Attentional Push-based augmented saliency with the standard saliency methods for both static and dynamic stimuli. The results show that each of the augmented saliency methods improves its corresponding saliency method and the average evaluation scores for the augmented saliency methods are significantly higher than the average scores of the standard saliency methods. For static stimuli, the most significant performance boost in *AUC* score is achieved by augmenting the AWS method (although the augmented Center model has the highest improvement, its *AUC* score is insignificant compared to the best performing method). The average performance boost over all of the augmented methods are 0.056, 0.42 and 0.11 for *AUC*, *NSS* and *CC*, respectively. It should be noted that the augmented saliency method is not only outperforming models that employ face and people detection such as Judd [36], it is also improving the prediction performance of data-driven methods such as the ensemble of Deep Networks (eDN) [23].

The performance improvements are more noticeable for the dynamic imagery. The average performance boosts for all of the augmented methods are 0.10, 1.19 and 0.18 for *AUC*, *NSS* and *CC*, respectively. The most significant performance boost in *AUC* score for the dynamic stimuli belongs to the augmented ContextAware model, which is more than 3 times larger than its improvement for static stimuli. This implies that the Attentional Push cues have more influence

Table 1: Average evaluation scores for the Attentional Push-based augmented saliency vs. ten best-performing saliency models on static and dynamic stimuli. The best performing method is shown in bold for each metric.

	AUC		NSS		CC	
	static	dynamic	static	dynamic	static	dynamic
AWS [13]	0.78	0.79	1.16	1.02	0.31	0.16
augmented AWS	0.85	0.91	1.66	<b>2.44</b>	0.44	<b>0.37</b>
BMS [22]	0.80	0.80	1.19	1.15	0.31	0.17
augmented BMS	0.85	0.90	1.63	2.30	0.43	0.35
Center [31]	0.61	0.75	0.47	0.99	0.13	0.15
augmented Center	0.77	0.90	1.20	2.26	0.32	0.35
ContextAware [11]	0.79	0.66	1.18	0.40	0.31	0.06
augmented ContextAware	0.85	0.88	1.61	2.10	0.43	0.31
eDN [23]	0.85	0.90	1.23	1.43	0.33	0.22
augmented eDN	<b>0.87</b>	<b>0.92</b>	1.58	2.21	0.42	0.34
FES [50]	0.82	0.83	1.49	0.97	0.39	0.15
augmented FES	0.85	0.89	<b>1.77</b>	2.16	<b>0.47</b>	0.33
GBVS [10]	0.81	0.85	1.31	1.36	0.35	0.21
augmented GBVS	0.85	0.90	1.61	2.29	0.43	0.35
IttiKoch2 [9]	0.79	0.80	1.17	1.04	0.31	0.16
augmented IttiKoch2	0.85	0.90	1.59	2.13	0.42	0.33
Judd [36]	0.84	0.87	1.30	1.34	0.35	0.21
augmented Judd	0.86	0.91	1.61	2.21	0.43	0.34
RARE [12]	0.80	0.75	1.25	0.54	0.33	0.08
augmented RARE	0.85	0.89	1.66	2.16	0.44	0.33
Average improvements	0.056	0.10	0.42	1.19	0.11	0.18

upon the viewers’ fixation in dynamic scenes, which could be explained by the observation that the viewers feel more immersed while watching dynamic scenes. Example saliency maps for some of the augmented and standard saliency methods are shown in Fig. 4.

To evaluate the effect of each Attentional Push cue in predicting the viewers’ fixation, we create five separate augmented saliency maps, each based on a single Attentional Push cue. We use the AWS model as the standard saliency method to compute the augmented saliency maps. Table 2 presents the average evaluation scores for the dynamic stimuli. Although the static Attentional Push cues seem to dominate most of the performance improvements, the dynamic Attentional Push cues have contribution in the performance improvements nonetheless. It should be noted that dynamic Attentional Push cues are not active in each frame and they require triggering event such as scene changes and changes in gaze direction. Given a saliency method augmented using only a dynamic Attentional Push cue, we can expect the average improvements over all the video frames to be small. Nevertheless, for a saliency map augmented using a combination of static and dynamic Attentional Push cues, the dynamic cues can make contributions in improving the performance on many video frames that would be missed by static

Table 2: Average evaluation scores of five separate augmented saliency maps, each based on a single Attentional Push cue for the dynamic stimuli.

	None	Static cues		Dynamic cues			All
		head/body	pose centerbias	head/body	pose Bounce	SceneChange	
AUC	0.79	0.87	0.82	0.80	0.80	0.81	0.91
NSS	1.02	1.53	1.32	1.28	1.19	1.13	2.44
CC	0.16	0.23	0.19	0.20	0.17	0.17	0.37

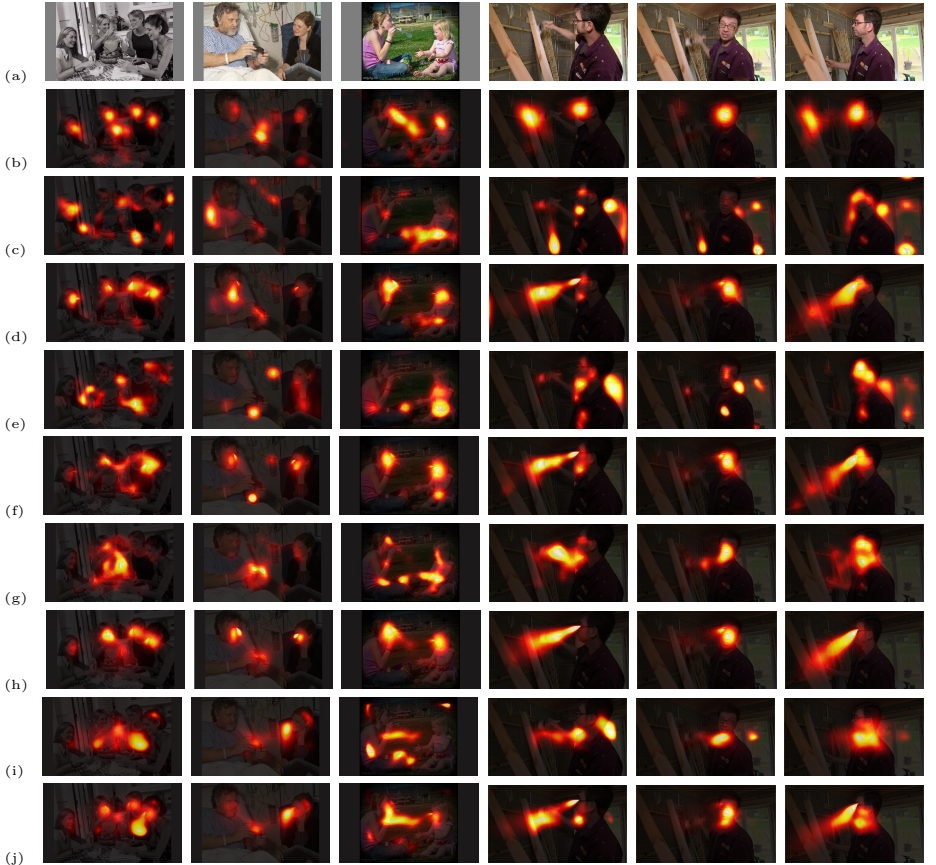


Fig. 4: Sample images and video frames from the CAT2000 [24] and the DIEM [28] dataset with overlaid ground truth, overlaid saliency maps and overlaid Attentional Push-based augmented saliency maps. (a) Original images, (b) overlaid ground truth, (c) overlaid AWS maps, (d) overlaid augmented AWS maps, (e) overlaid BMS maps, (f) overlaid augmented BMS maps, (g) overlaid eDN maps, (h) overlaid augmented eDN maps, (i) overlaid FES maps, (j) overlaid augmented FES maps. Augmented saliency methods alter the standard saliency maps to be more consistent with the ground truth.

cues. It can be seen in Table 2 that the combination of static and dynamic cues clearly outperforms static cues.

We examined the cases in which the prediction performance of the augmented saliency map is lower than the saliency map in static stimuli. For each static stimulus, we consider images for which at least two of the three evaluation scores display degraded performance. There are twelve such images in total, with two of them showing degraded performance consistently in all evaluation metrics. Both of these images contain crowded scenes, in which the actors are looking in many different directions. The reason for the degraded performance lies in the fact that the scene actors do not share the same loci of attention and therefore, the Attentional Push cues arising from their gaze directions compete with one another in pushing the viewers' attention. This situation leads to an inconsistent increase in the saliency values of many image regions that are not foci of actors' attention, which would lead to a degraded prediction performance for the augmented saliency method.

## 5 Conclusion

We presented an attention modeling scheme which combines Attentional Push cues, i.e. the power of image regions to direct and manipulate the attention allocation of the viewer, with standard saliency models, which generally concentrate on analyzing image regions for their power to pull attention. Our methodology significantly outperforms saliency methods in predicting the viewers' fixations on both static and dynamic stimuli. Our results showed that by employing Attentional Push cues, the augmented saliency maps can challenge the state of the art in saliency models.

## References

1. Tsotsos, J.K.: Analyzing vision at the complexity level. *Behavioral and Brain Sciences* **13** (1990) 423–445
2. Walther, D., Koch, C.: Modeling attention to salient proto-objects. *Neural Networks* **19**(9) (2006) 1395 – 1407
3. Benfold, B., Reid, I.: Guiding visual surveillance by tracking human attention. In: *Proceedings of the 20th British Machine Vision Conference*. (2009) 1–11
4. Ferreira, J., Dias, J.: Attentional mechanisms for socially interactive robots- a survey. *Autonomous Mental Development, IEEE Transactions on* **6**(2) (June 2014) 110–125
5. Rosenholtz, R., Dorai, A., Freeman, R.: Do predictions of visual perception aid design? *ACM Trans. Appl. Percept.* **8**(2) (February 2011) 12:1–12:20
6. Treisman, A.M., Gelade, G.: A feature-integration theory of attention. *Cognitive Psychology* **12**(1) (1980) 97–136
7. Koch, C., Ullman, S.: Shifts in selective visual attention: Towards the underlying neural circuitry. *Human Neurobiology* **4**(4) (1985) 219–227
8. Clark, J.J., Ferrier, N.J.: Modal control of an attentive vision system. In: *Computer Vision., Second International Conference on*. (Dec 1988) 514–523

9. Itti, L., Koch, C., Niebur, E.: A model of saliency-based visual attention for rapid scene analysis. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **20**(11) (1998) 1254–1259
10. Harel, J., Koch, C., Perona, P.: Graph-based visual saliency. In: *NIPS*. (2007) 545–52
11. Goferman, S., Zelnik-Manor, L., Tal, A.: Context-aware saliency detection. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **34**(10) (2012) 1915–1926
12. Riche, N., Mancas, M., Duvinage, M., Mibulumukini, M., Gosselin, B., Dutoit, T.: Rare2012: A multi-scale rarity-based saliency detection with its comparative statistical analysis. *Signal Processing: Image Communication* **28**(6) (2013) 642–658
13. Garcia-Diaz, A., Leborn, V., Fdez-Vidal, X.R., Pardo, X.M.: Corrections to: On the relationship between optical variability, visual saliency, and eye fixations: A computational approach. *Journal of Vision* **12**(7) (2012) 13
14. Zhang, J., Sclaroff, S.: Saliency detection: A boolean map approach. In: *Computer Vision (ICCV), 2013 IEEE International Conference on*. (2013) 153–160
15. Borji, A., Itti, L.: State-of-the-art in visual attention modeling. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(1) (2013) 185–207
16. Haji-Abolhassani, A., Clark, J.J.: An inverse yarbus process: Predicting observers task from eye movement patterns. *Vision Research* **103** (2014) 127–142
17. Borji, A., Itti, L.: Defending yarbus: Eye movements reveal observers’ task. *Journal of Vision* **14**(3) (2014) 29
18. Kanan, C., Ray, N.A., Bseiso, D.N.F., Hsiao, J.H., Cottrell, G.W.: Predicting an observer’s task using multi-fixation pattern analysis. In: *Proceedings of the Symposium on Eye Tracking Research and Applications. ETRA ’14*, New York, NY, USA (2014) 287–290
19. Torralba, A., Castelano, M.S., Oliva, A., Henderson, J.M.: Contextual guidance of eye movements and attention in real-world scenes: the role of global features in object search. *Psychological Review* **113**(4) (2006) 766–786
20. Tatler, B.W.: The central fixation bias in scene viewing: Selecting an optimal viewing position independently of motor biases and image feature distributions. *Journal of Vision* **7**(4) (2007) 1–17
21. Underwood, G., Foulsham, T., Humphrey, K.: Tsaliency and scan patterns in the inspection of real-world scenes: Eye movements during encoding and recognition. *Visual Cognition* **17** (2009) 812–834
22. Zhang, J., Sclaroff, S.: Exploiting surroundedness for saliency detection: A boolean map approach. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2015)
23. Vig, E., Dorr, M., Cox, D.: Large-scale optimization of hierarchical features for saliency prediction in natural images. In: *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*. (June 2014) 2798–2805
24. Borji, A., Itti, L.: CAT2000: A large scale fixation dataset for boosting saliency research. *CVPR 2015 workshop on "Future of Datasets"* (2015)
25. Borji, A., Sihite, D., Itti, L.: Quantitative analysis of human-model agreement in visual saliency modeling: A comparative study. *IEEE Trans. Image Processing* **22**(1) (2013)
26. Birmingham, E., Bischof, W.F., Kingstone, A.: Saliency does not account for fixations to eyes within social scenes. *Vision Research* **49**(24) (2009) 2992 – 3000
27. Borji, A., Parks, D., Itti, L.: Complementary effects of gaze direction and early saliency in guiding fixations during free viewing. *Journal of Vision* **14**(13) (2014) 1–32

28. Mital, P.K., Smith, T.J., Hill, R., Henderson, J.M.: Clustering of gaze during dynamic scene viewing is predicted by motion. *Cognitive Computation* **3**(1) (2011) 5–24
29. Kaplan, F., Hafner, V.V.: The challenges of joint attention. *Interaction Studies* **7**(2) (2006) 135–169
30. Smith, K., Ba, S., Odobez, J., Gatica-Perez, D.: Tracking the visual focus of attention for a varying number of wandering people. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **30**(7) (July 2008) 1212–1229
31. Bylinskii, Z., Judd, T., Borji, A., Itti, L., Durand, F., Oliva, A., Torralba, A.: Mit saliency benchmark. <http://saliency.mit.edu>
32. Parks, D., Borji, A., Itti, L.: Augmented saliency model using automatic 3D head pose detection and learned gaze following in natural scenes. *Vision Research* **116**, Part B (2015) 113 – 126
33. Castelhamo, M.S., Wieth, M., Henderson, J.M.: I see what you see: Eye movements in real-world scenes are affected by perceived direction of gaze. In: *Attention in Cognitive Systems. Theories and Systems from an Interdisciplinary Viewpoint*. Springer Berlin Heidelberg (2007) 251–262
34. Ricciardelli, P., Bricolo, E., Aglioti, S.M., Chelazzi, L.: My eyes want to look where your eyes are looking: Exploring the tendency to imitate another individuals gaze. *Neuroreport* **13**(17) (2002) 2259–2264
35. Tseng, P.H., Carmi, R., Cameron, I.G.M., Munoz, D.P., Itti, L.: Quantifying center bias of observers in free viewing of dynamic natural scenes. *Journal of Vision* **9**(7) (2009) 4
36. Judd, T., Ehinger, K., Durand, F., Torralba, A.: Learning to predict where humans look. In: *Computer Vision, 2009 IEEE 12th International Conference on*. (2009) 2106–2113
37. Smith, T.J.: The attentional theory of cinematic continuity. *Projections* **6**(1) (2012) 1–27
38. Hoffman, M.W., Grimes, D.B., Shon, A.P., Rao, R.P.: A probabilistic model of gaze imitation and shared attention. *Neural Networks* **19** (2006) 299–310
39. Moon, A., Troniak, D.M., Gleeson, B., Pan, M.K., Zheng, M., Blumer, B.A., MacLean, K., Croft, E.A.: Meet me where i'm gazing: How shared attention gaze affects human-robot handover timing. In: *Proceedings of the 2014 ACM/IEEE International Conference on Human-robot Interaction*. (2014) 334–341
40. Triesch, J., Teuscher, C., Carlson, E.: Gaze following: Why (not) learn it. *Developmental Science* **9**(2) (2006) 125–147
41. Marat, S., Rahman, A., Pellerin, D., Guyader, N., Houzet, D.: Improving visual saliency by adding face feature map and center bias. *Cognitive Computation* **5**(1) (2013) 63–75
42. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. Volume 1. (June 2005) 886–893 vol. 1
43. Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features. In: *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*. Volume 1. (2001) I–511–I–518 vol.1
44. Zhu, X., Ramanan, D.: Face detection, pose estimation and landmark localization in the wild. In: *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. (2012)
45. Kalal, Z., Mikolajczyk, K., Matas, J.: Tracking-learning-detection. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **34**(7) (2012) 1409–1422



46. Asthana, A., Zafeiriou, S., Cheng, S., Pantic, M.: Incremental face alignment in the wild. In: Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on. (2014) 1859–1866
47. Maji, S., Bourdev, L., Malik, J.: Action recognition from a distributed representation of pose and appearance. In: Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on. (2011)
48. Kim, Y.M., Choi, S.W., Lee, S.W.: Fast scene change detection using direct feature extraction from mpeg compressed videos. In: Pattern Recognition, 2000. Proceedings. 15th International Conference on. Volume 3. (2000) 174–177 vol.3
49. Ramanathan, S., Katti, H., Sebe, N., Kankanhalli, M., Chua, T.S.: An eye fixation database for saliency detection in images. In: ECCV 2010, Crete, Greece (2010)
50. Tavakoli, H.R., Rahtu, E., Heikkilä, J.: Fast and efficient saliency detection using sparse sampling and kernel density estimation. In: Proceedings of the 17th Scandinavian Conference on Image Analysis. SCIA'11, Springer-Verlag (2011) 666–675